# Genome project management resources at the National Agricultural Library

Monica Poelchau/Chris Childers, USDA-NAL

Monica Muñoz-Torres, BBOP/LBNL

Aphid Genomics Symposium. September 25th, 2016

USDA

# So you have a genome project. Where will you store your data?

- Make your data available through NCBI (or other INSDC organizations).

- To make your data even more useful for your community, consider also making it available in a taxon-specific repository.

- Advantages:
  - Greater visibility for your dataset
  - Easier to find data for comparative analyses
  - Value-added tools for searching and browsing, analysis
  - Curation tools to improve annotation quality

USDA

# The i5k Workspace@NAL

## Our focus:

- We support **any** 'orphaned' arthropod genome project:
  - Genome assembly needs to be in GenBank/ENA/DDBJ
  - Data should be open access (no private repositories)
- We enable and support **community curation.**

## Our background:

- Originally set up to support genomes sequenced as part of the i5k initiative
- To learn more about the i5k initiative, visit **Booth #320**

# The i5k Workspace@NAL

- ## Resources:
  - Organism landing page (Tripal software)
  - Gene pages for official gene sets
  - Tutorials

- ## Tools:
  - BLAST, HMMER, Clustal
  - JBrowse genome browser
  - Apollo curation software

- ## Support:
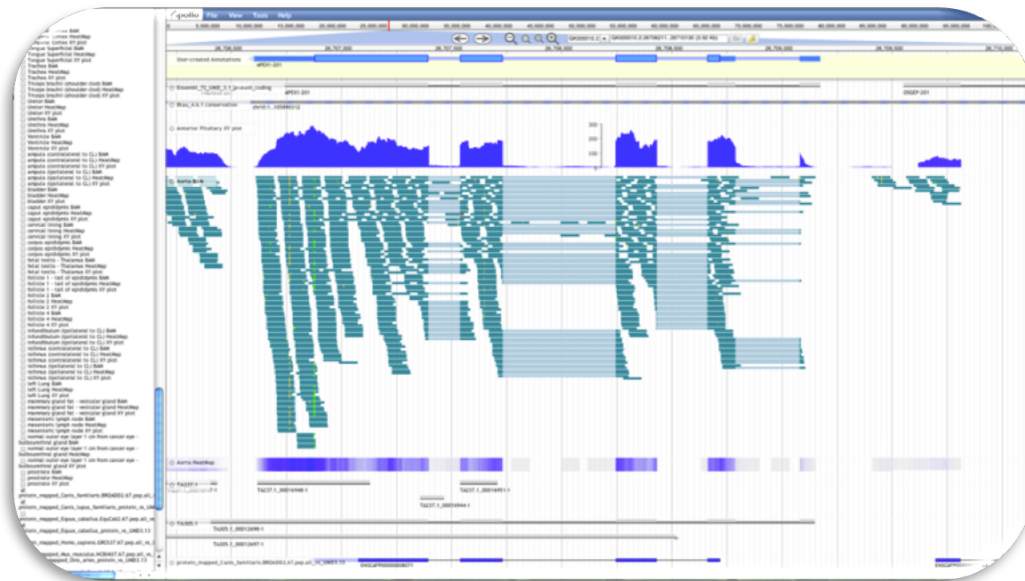  - Semi-automated QC of manual annotations
  - OGS generation pipeline

# i5k Workspace data – 53 species and counting

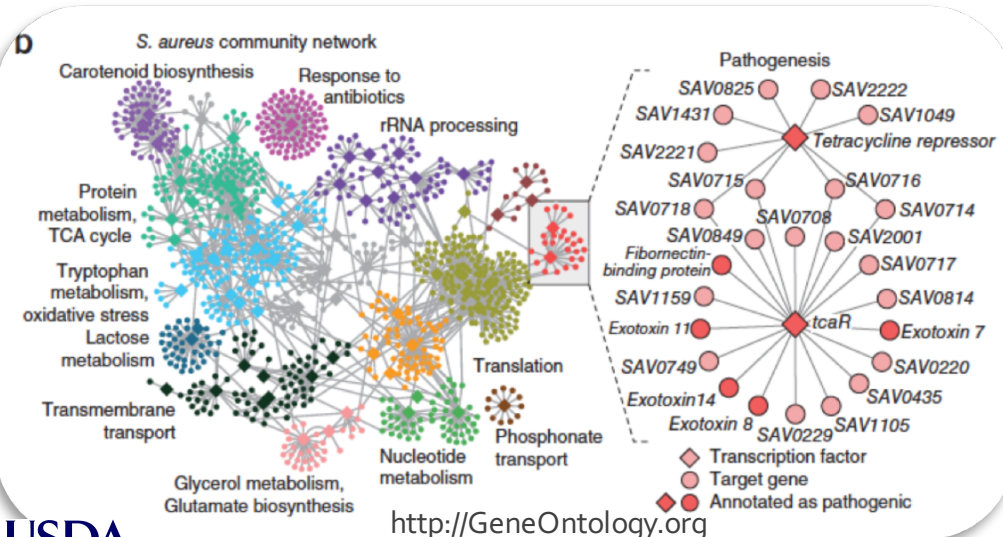| Order | Quantity | Order | Quantity |
|---|---|---|---|
| Amphipoda | 1 | Harpacticoida | 1 |
| Araneae | 3 | Hemiptera | 7 |
| Blattodea | 1 | Hymenoptera | 13 |
| Calanoida | 1 | Lepidoptera | 2 |
| Coleoptera | 5 | Odonata | 1 |
| Diplura | 1 | Scorpiones | 1 |
| Diptera | 13 | Thysanoptera | 1 |
| Ephemeroptera | 1 | Trichoptera | 1 |

- Many other datasets mapped to, or predicted from each genome assembly (gene predictions, transcriptomes, RNA-Seq, etc.)

USDA

# Curation



**Identifies** elements that best represent the underlying biology & **eliminates** elements that reflect systemic errors of automated analyses.

**Assigns function** through comparative analysis of similar genome elements from closely related species using literature, databases, and experimental data.



http://GeneOntology.org

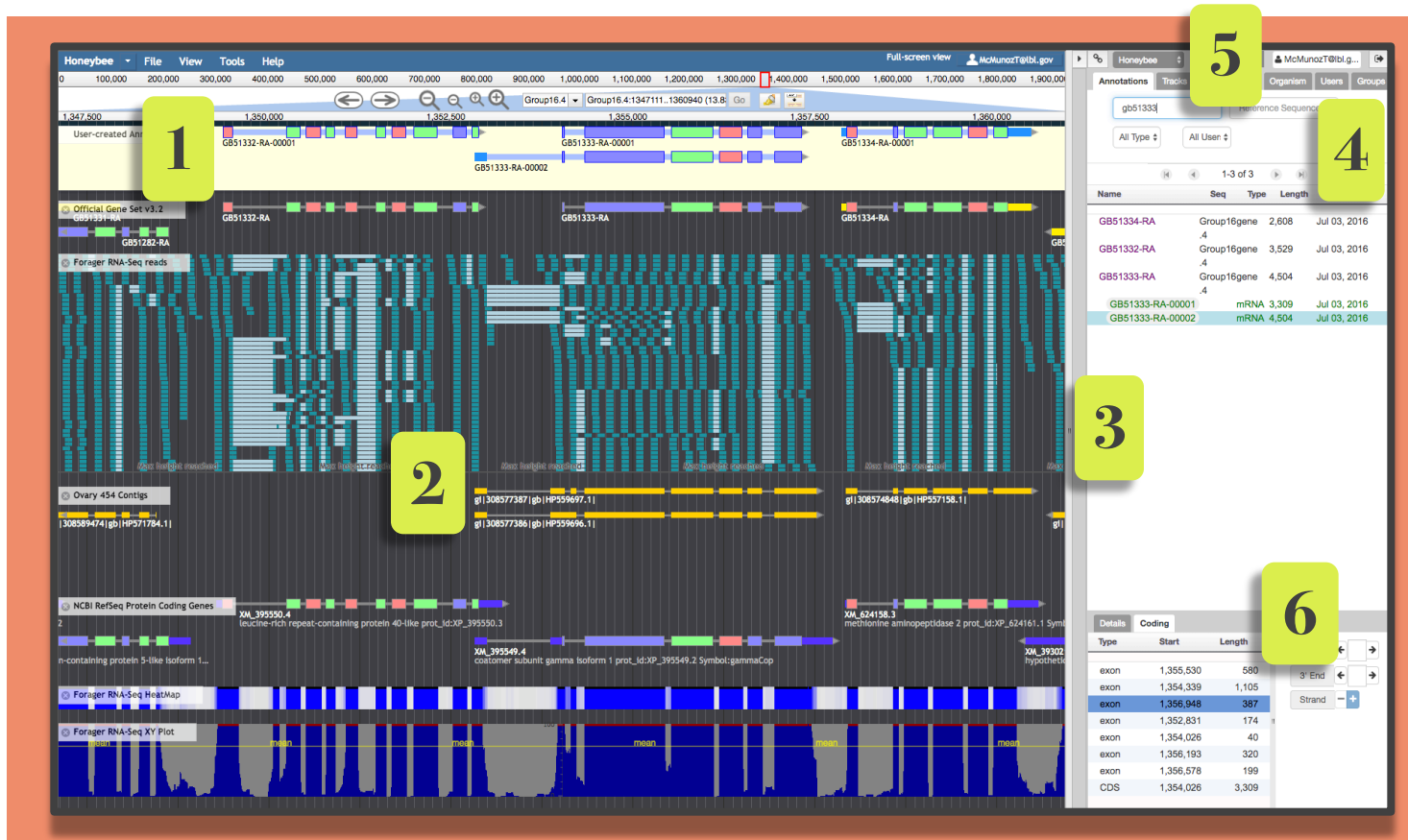# Community curation at the i5k Workspace

- Why curate?
  - Verify quality of automated gene predictions
  - Improve gene models for specific analyses

- **Our community:** Over 400 registered annotators have curated more than10,000 gene models using the Apollo genome annotation editor.

# Community curation at the i5k Workspace

Our support for community curation includes:

- Access to a large community of curators

- Tutorials, guidelines, webinars

- Registration mechanism for new annotators

- One-on-one support

- Software to evaluate changes between curated and original annotations (Chien-Yueh Lee, https://github.com/chienyuehlee/gff-cmp-cat)

# Apollo: Collaborative, instantaneous, web-based



1. User-created Annotations
2. Evidence Tracks: Experimental data, alignments
3. Annotator Panel: Removable dock
4. Tabs for searching, editing, and exporting data
5. Switch between organisms and sequences
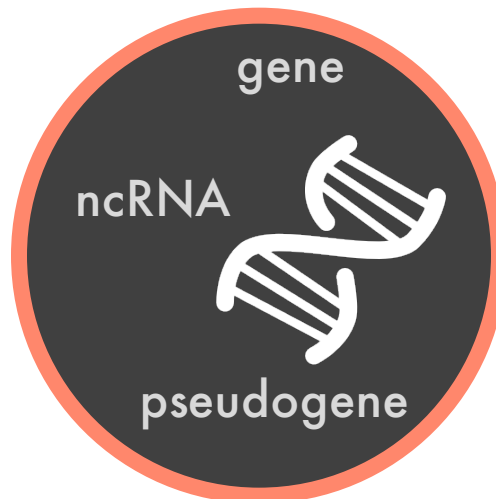6. Visualize and edit annotation details

# Apollo: Latest Features



Annotate multiple organisms per server

Change type of genomic element to annotate

gene

ncRNA

pseudogene

Galaxy / Apollo Integration

Export and update a Chado database

# QC and OGS pipeline

- QC program corrects common formatting errors from the curation process

- OGS generation program merges curated models with one designated gene set using curator-supplied information
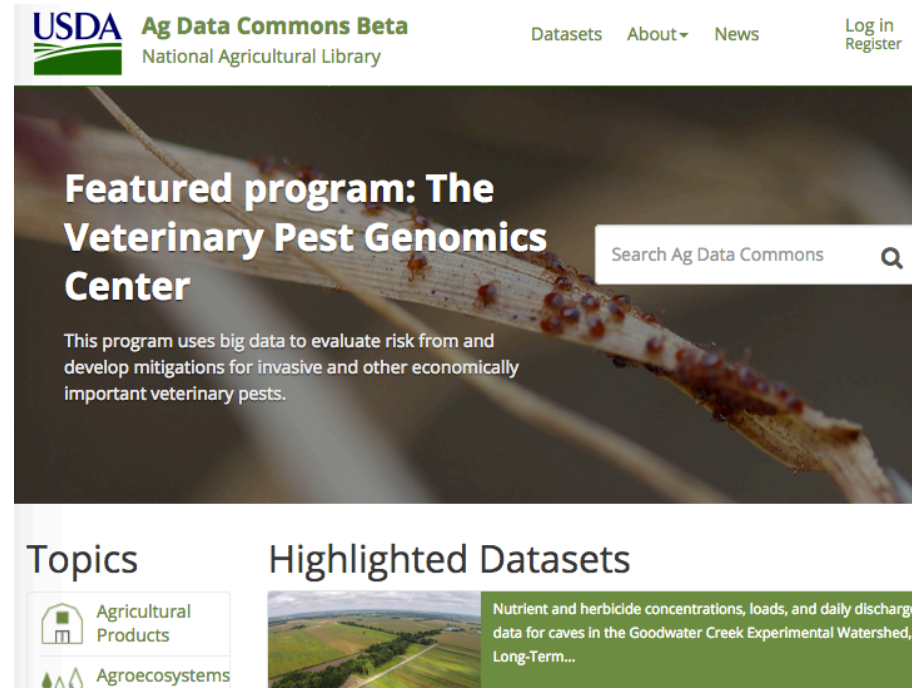
- Still in development, already 4 OGS's produced (Mei-Ju Chen)

Apollo output → Error checking  Curator fixes → Merge with one designated gene set → Official Gene Set

# Genome already hosted elsewhere?

- You can also use our tools to query the datasets that we host.

# Other resources at the NAL: The Ag Data Commons

- Hosts any dataset funded by the USDA
- Landing page
- Citable DOI
- https://data.nal.usda.gov/
- Nine i5k datasets already available

# Need more information?

## i5k Workspace@NAL:
- https://i5k.nal.usda.gov/
- https://github.com/NAL-i5K/

## The i5k initiative:
- New website: http://i5k.github.io/
- New webinar series coming soon!

## Apollo:
- http://GenomeArchitect.org/

## Ag Data Commons:
- https://data.nal.usda.gov/

Learn more about i5k Workspace@NAL at Poster D3385 on Tuesday

Visit us this week at Booth 320; 1-5 PM

USDA

# Acknowledgements

**The NAL Team**

- Vijaya Tsavatapalli
- Gary Moore
- Susan McCarthy
- Yu-yu Lin
- Mei-Ju Chen

**Workspace alumni**

- Chien-Yueh Lee
- Han Lin
- Jun-Wei Lin

**i5k Workspace@NAL advisory committee**

- Jay Evans
- Kevin Hackett
- Simon Liu
- Ursula Pieper

- i5k Coordinating Committee
- i5k Pilot Project
- Apollo & JBrowse Development Teams
- GMOD/Tripal community
- **All of our users and contributors!**